

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2002-041078

(43)Date of publication of application : 08.02.2002

(51)Int.Cl.

G10L 15/06  
G10L 15/10  
G10L 15/00  
G10L 15/28

(21)Application number : 2000-220576

(71)Applicant : SHARP CORP

(22)Date of filing : 21.07.2000

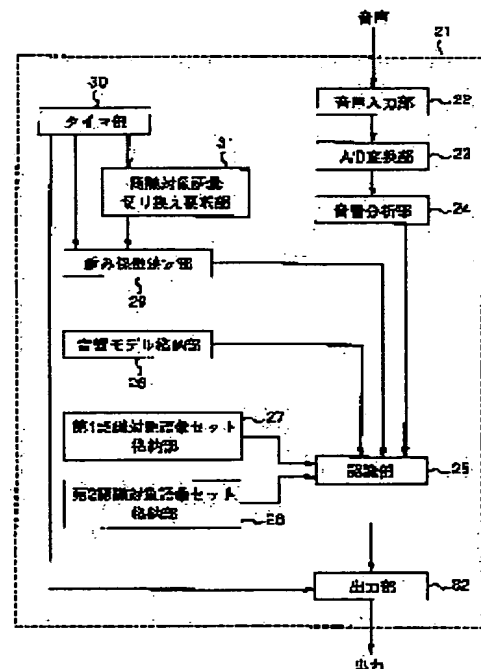
(72)Inventor : HONDA KAZUMASA  
TSURUTA AKIRA  
KANZA HIROYUKI

## (54) VOICE RECOGNITION EQUIPMENT, VOICE RECOGNITION METHOD AND PROGRAM RECORDING MEDIUM

## (57)Abstract:

PROBLEM TO BE SOLVED: To obtain high recognition accuracy even when the recognition object vocabulary is automatically changed.

SOLUTION: A recognition means 25 calculates the likelihood P of the vocabulary which constitutes the recognition object vocabulary sets A, B stored in the first and second recognition object vocabulary storing parts 27, 28 using the acoustics model in an acoustics model storing part 26. The change of the recognition object vocabulary sets accompanied with the display contents in an output part 32 is performed at the time by changing the values of weights  $w_1$ ,  $w_2$  which multiply to P between '1' and the appointed value 'a' near 0 (zero) proportioned the passing time from the requested change time  $t_0$ . As the result, even when the speaker misses the utterance chance and the recognition object vocabulary is changed automatically, high recognition result can be obtained if the speaker utters the recognition object vocabulary before changing because the calculations of the likelihood  $w \times P$  are also performed.



## LEGAL STATUS

[Date of request for examination] 26.07.2002

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number] 3563018

[Date of registration] 11.06.2004

[Number of appeal against examiner's decision  
of rejection]

[Date of requesting appeal against examiner's  
decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開2002-41078

(P2002-41078A)

(43) 公開日 平成14年2月8日 (2002.2.8)

(51) Int.Cl. <sup>7</sup>	識別記号	F I	テーマコード(参考)
G 1 0 L 15/06		G 1 0 L 3/00	5 2 1 V 5 D 0 1 5
15/10			5 3 1 G
15/00			5 5 1 P
15/28			

審査請求 未請求 請求項の数 6 O L (全 11 頁)

(21) 出願番号 特願2000-220576(P2000-220576)

(22) 出願日 平成12年7月21日 (2000.7.21)

(71) 出願人 000005049

シャープ株式会社

大阪府大阪市阿倍野区長池町22番22号

(72) 発明者 本田 和正

大阪府大阪市阿倍野区長池町22番22号 シ

ャープ株式会社内

(72) 発明者 鶴田 彰

大阪府大阪市阿倍野区長池町22番22号 シ

ャープ株式会社内

(74) 代理人 100062144

弁理士 青山 葆 (外1名)

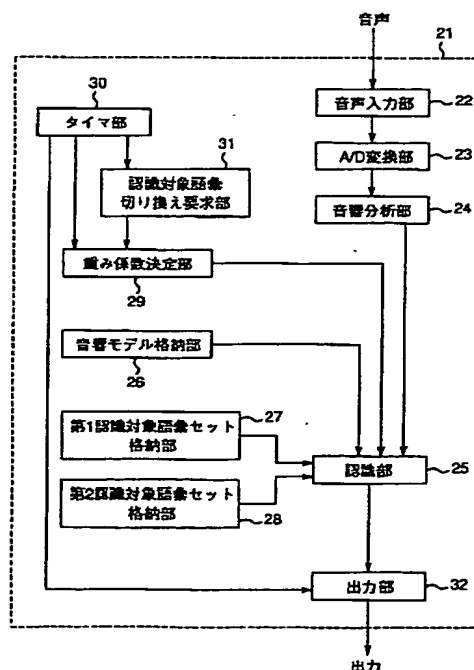
最終頁に続く

(54) 【発明の名称】 音声認識装置および音声認識方法、並びに、プログラム記録媒体

(57) 【要約】

【課題】 自動的に認識対象語彙を切り換える場合でも高認識精度を得る。

【解決手段】 認識部25は、音響モデル格納部26の音響モデルを用いて第1,第2認識対象語彙セット格納部27に格納された認識対象語彙セットA,Bを構成する単語の尤度Pを算出する。その際に、出力部32の表示内容の切り換えに伴う認識対象語彙セットの切り換えは、切り換え前後の認識対象語彙セットを構成する単語の尤度Pに掛ける重み $w_1, w_2$ の値を切り換え要求時刻 $t$ からの経過時間に比例して「1」と0近傍の所定値「a」との間で切り換えることで行う。その結果、話者が認識対象語彙の発声の機会を逸し、且つ、自動的に認識対象語彙が切り換えられても、切り換え前の認識対象語彙セットの単語に関する尤度 $w \cdot P$ の計算も行われるために、話者が切り換え前の認識対象語彙で発声しても高い認識結果が得られる。



## 【特許請求の範囲】

【請求項 1】 入力された音声を認識する認識部と、この認識部の認識結果を含む情報を出力する出力部と、上記認識時に用いられる認識対象語彙が格納された認識語彙格納部と、タイマ部と、このタイマ部からの時刻信号に基づいて上記認識対象語彙の切り換えを要求する認識対象語彙切り換え要求部を有する音声認識装置において、

上記出力部は、複数の出力内容を切り換え出力するようになっており、

上記認識対象語彙は、上記出力部の出力内容に対応した認識対象語の集合でなる複数の認識対象語彙セットに分類され、上記認識対象語彙の切り換えは上記認識対象語彙セットの単位で行われるようになっており、

上記タイマ部からの時刻信号に基づいて、上記各認識対象語彙セット用の重みを決定する重み決定部を備えて、上記認識部は、上記全認識対象語彙セットおよび上記決定された各重みを用いて、入力音声を認識するようになっていることを特徴とする音声認識装置。

【請求項 2】 請求項 1 に記載の音声認識装置において、

上記重み決定部は、上記認識対象語彙切り換え要求部によって認識対象語彙の切り換えが要求されてから重み決定までの経過時間に応じて、切り換え前の認識対象語彙セット用の重みを低下させる一方、切り換え後の認識対象語彙セット用の重みを上昇させるようになっていることを特徴とする音声認識装置。

【請求項 3】 請求項 1 あるいは請求項 2 に記載の音声認識装置において、

上記認識部は、上記全認識対象語彙セットを構成する各語の尤度を算出し、各語の尤度の値に各語が属する認識対象語彙セット用の重みを掛け、その値が最も高い語を認識結果とするようになっていることを特徴とする音声認識装置。

【請求項 4】 請求項 2 に記載の音声認識装置において、

上記出力部は、上記認識対象語彙切り換え要求部からの認識対象語彙切り換え要求がなされた時点に出力している出力内容に対応する認識対象語彙セット用の重みの値と、次に出力すべき出力内容に対応する認識対象語彙セット用の重みの値との差が所定値未満になると、上記出力内容を切り換えるようになっていることを特徴とする音声認識装置。

【請求項 5】 入力された音声を認識対象語彙を用いて認識して認識結果を出力するに際して、タイマ部からの時刻信号に基づいて上記認識対象語彙の切り換えを自動的に行う音声認識方法において、

複数の出力内容を出力部に切り換え出力し、

上記各出力内容に対応した認識対象語の集合でなる複数の認識対象語彙セットの単位で、上記認識対象語彙の切

り換えを行い、

上記タイマ部からの時刻信号に基づいて、上記各認識対象語彙セット用の重みを決定し、

上記全認識対象語彙セットおよび上記決定された各重みを用いて、上記入力音声の認識を行うことを特徴とする音声認識方法。

【請求項 6】 コンピュータを、

請求項 1 における認識部、出力部、タイマ部、認識対象語彙切り換え要求部および重み決定部として機能させる音声認識処理プログラムが記録されたことを特徴とするコンピュータ読出し可能なプログラム記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】この発明は、コンピュータや携帯情報端末に搭載されて人間の発声による音声を認識する音声認識装置および音声認識方法、並びに、音声認識処理プログラムを記録したプログラム記録媒体に関する。

【0002】

【従来の技術】音声認識装置において、認識精度を高めるために、必要に応じて認識対象語彙を切り換えるという認識方法がある。このような認識方法を用いた音声認識装置の応用例として、パーソナルコンピュータや日本語ワードプロセッサ等の表示装置を有する機器において、表示装置を用いたメニュー表示による機器の操作ガイドを、音声認識を用いて行うことが考えられる。

【0003】上述のような操作ガイドによれば、操作方法や操作による効果の表示を画面で確認しながら操作を学ぶことができる。そして、上記表示装置の画面が狭い等、上記表示装置からの情報量が少ない場合には、複数の機器操作に関する操作ガイドの表示を時間の経過と共に自動的に切り換える場合がある。このような操作ガイドに音声を用いれば、利用者にとって分かり易く、且つ、操作ボタンの数を減らして操作を簡単にすることができる。その場合、複数の機器操作に関する操作ガイドの表示の切り換えと共に認識対象語彙を切り換えれば、認識対象語彙を少なくすることができるので高い認識精度を得ることができる。

【0004】このような認識対象語彙を切り換える認識方法の応用においては、切り換え表示する各メニューに関連のある認識対象語彙のセットを、メニュー数分だけ複数記憶しておく。そして、利用者の操作や時間の経過等によるメニュー表示の切り換えに同期して認識対象語彙を切り換えることによって、夫々のメニューにおいては必要最小限の語彙を対象に認識処理を行うことができ、認識精度を向上させることができるのである。その場合、時間の経過と共にメニュー表示を自動的に切り換える際には、機器が自動的に認識対象語彙をも切り換えることになる。

【0005】以下、上記認識対象語彙の切り換えが可能

な音声認識装置について説明する。図4は、上記認識対象語彙の切り換えが可能な音声認識装置の一例を示すブロック図である。ここで、本音声認識装置1は、認識対象語彙の切り換えおよび出力部13による表示内容の切り換えは、所定時間毎に音声認識装置1自身が自動的に行うものとする。音声認識装置1は、A/D(アナログ/デジタル)変換部2、音響分析部3、認識部4、音響モデル格納部5、認識対象語彙格納・判定部6、現認識対象語彙識別子記憶部7、タイマ部8、認識対象語彙切り換え要求部9、認識対象語彙切り換え要求時刻記憶部10、音声検出部11、音声時刻記憶部12、出力部13で構成される。

【0006】話者によって上記音声認識装置1に入力された音声は、A/D変換部2に送出されてデジタル化される。そして、このデジタル化された音声波形は、音響分析部3で、20msec〜40msecの区間毎に比較的短時間の時間窓を掛けると共に、8msec〜16msec毎に上記時間窓をシフトしていく短時間スペクトル分析の手法によって分析される。上記時間窓によって切り出された音声波形は、切り出し時の時間長を有するフレームと呼ばれる単位の特徴ベクトルの時系列に変換される。ここで、上記特徴ベクトルは、その時刻における音声スペクトルの特徴量を抽出したもので、通常は10次元〜100次元であり、LPC(線形予測分析)メルケプストラム係数等が広く用いられている。こうして変換された特徴ベクトルは認識部4に送出されると共に、音声入力開始を検出する音声検出部11にも出力される。そうすると、音声時刻記憶部12は、音声検出部11からの音声入力開始信号とタイマ部8からの時刻信号とに基づいて、音声入力開始時刻を検出して記憶する。

【0007】上記音響モデル格納部5には、認識単位毎に用意されたHMM(隠れマルコフモデル)が用意されている。ここで、上記認識単位としては、音素や単語が広く用いられている。また、HMMとは、複数の状態を有する非決定性確率有限オートマトンであり、非定常信号源を定常信号源の連結で表す統計的信号源モデルである。尚、出力確率や遷移確率等のパラメータは、対応する学習音声を与えてバウム・ウェルチアルゴリズムと呼ばれるアルゴリズム等によって予め学習されている。以下、音響モデル格納部5には、認識単位が音素であるHMMが記憶されているものとする。

【0008】上記認識対象語彙の切り換えの動作は、特開平6-337695号公報に開示されている方法を適用する。上記認識対象語彙として、認識対象語彙セットAと認識対象語彙セットBとがあり、現時点においては認識対象語彙識別子記憶部7には認識対象語彙セットAの識別子が記憶されているものとする。また、出力部13は、認識対象語彙セットAに対応する表示内容を表示しているものとする。

【0009】この状態で、所定時間が経過すると、タイ

マ部8から認識対象語彙切り換え要求部9及び出力部13に対して通知がなされる。そうすると、出力部13は、表示内容を認識対象語彙セットBに対応する表示内容に変更する。また、認識対象語彙切り換え要求部9から切り換えが要求され、その要求時刻が認識対象語彙切り換え要求時刻記憶部10に記憶される。そして、認識対象語彙格納・判定部6によって、認識対象語彙切り換え要求時刻記憶部10に記憶されている要求時刻Tcと音声時刻記憶部12に記憶されている音声入力開始時刻Tsとが比較され、音声入力開始時刻Tsが要求時刻Tcよりも後である場合には、認識対象語彙の切り換えが要求された後に発声が行われたのであるから適切な認識対象語彙セットは認識対象語彙セットBであると判定される。それ以外は、認識対象語彙セットAであると判定される。そして、該当する認識対象語彙セットの識別子で現認識対象語彙識別子記憶部7の記憶内容を更新するのである。

【0010】こうして、適切な認識対象語彙セットの判定が終了すると、認識部4は、音響分析部3で得られた特徴ベクトルと、現認識対象語彙識別子記憶部7に記憶されている識別子に対応して認識対象語彙格納・判定部6から出力される何れかの認識対象語彙セットを構成する各単語の音素列と、音響モデル格納部5に格納されているHMMを用いて、以下のようにして音声認識を行う。

【0011】すなわち、先ず、上記認識対象語彙に含まれる各単語のHMMを求める。具体的には、音響モデル格納部5に記憶されている各音素のHMMを、認識対象語彙セットを構成している各単語の音素列に対応させて結合するのである。

【0012】次に、夫々の単語のHMMについて、音響分析部3からの特徴ベクトルを用いて生起確率を求める。HMMによる音声認識においては、音声は初期状態から最終状態までの状態遷移の間にHMMから出力されるシンボルの時系列として表される。そこで、初期状態の確率を任意の値に定め、順次状態遷移毎に出力確率および遷移確率を掛けていくことによって、発声とそのモデルM(単語のHMM)から発生される確率を求めることができる。逆に、発声を観測した場合に、その発声があるモデルMから発生したと仮定すると、そのモデルMからの発生の確率が計算できることになる。

【0013】以下、上記認識部4における認識アルゴリズムについて詳細に説明する。認識部4は、音響分析部3によって得られた特徴ベクトルの時系列を入力とし、認識対象語彙格納・判定部6からの認識対象語彙に含まれる総ての単語のHMMに関してその生起確率を求め、最も高い生起確率を呈するHMMの単語を認識結果とする。すなわち、 $t(=1, 2, \dots, l)$ をフレーム番号として、特徴ベクトルの時系列で表現された入力の系列を、

$X = x_{\dots 1}, x_{\dots 2}, x_{\dots 3}, \dots, x_{\dots t}, \dots, x_{\dots l}$

とする。尚、「 $x_{vec} i$ 」は多次元のベクトルである。以下、ベクトル $x$ を「 $x_{vec}$ 」と表記する。さらに、モデル $M$ の初期状態の集合を $S$ とし、最終状態の集合を $F$ とする。また、「 $i, j$ 」を状態番号として、 $j$ 番目の状態の遷移系列を

$$Q = q_0^j, q_1^j, q_2^j, \dots, q_t^j, \dots, q_1^j$$

と表す。上式において、「 $q_t^j$ 」は、 $t$ 番目のフレームの入力記号 $x_{vec} t$ によって遷移した状態を表す。ここで、 $q_0^j \in S$ であり、 $q_1^j \in F$ である。更に、初期状態の初期確率を $\pi_i : \sum_i \pi_i = 1$ で表し、状態 $q_i$ から状態 $q_j$ への遷移確率を $a_{ij}$ とし、そのときに $x_{vec} i$ が出力される出力確率を $b_{ij}(x_{vec} i)$ とすると、入力系列の生起確率(尤度) $P(X|M)$ は、

$$P(X|M) = \sum_{all Q} \pi_0^j \prod_{i=1}^T a_{i-1,i}^j \cdot b_{i-1,i}^j(x_{vec} i)$$

で表される。この生起確率(尤度) $P(X|M)$ の演算を、認識対象語彙に含まれる全単語に対応するHMMについて計算し、最も高い生起確率(尤度) $P$ を呈するHMMに対応する単語を認識結果として出力部13に出力して表示するのである。

【0014】

【発明が解決しようとする課題】しかしながら、上記従来の特開平6-337695号公報に開示された認識対象語彙切り換え動作を適用した音声認識装置には、以下のような問題がある。すなわち、上述したように、特開平6-337695号公報に開示された認識対象語彙切り換え動作においては、音声入力開始時刻 $T_s$ が認識対象語彙切り換え要求時刻 $T_c$ よりも後である場合に認識対象語彙のセットを切り換えるようにしている。この方法は、話者の操作によって認識対象語彙切り換え要求がなされる場合には、必ず認識対象語彙の切り換え要求がなされた後に発声が行われるために有効である。

【0015】ところが、図4に示す音声認識装置のように、時間の経過と共に自動的に認識対象語彙が切り換る音声認識装置の場合には、認識対象語彙の切り換えは、話者の意識とは全く関係なく行われる。したがって、何らかの理由で話者が認識対象語彙の発声の機会を逸してしまい、且つ、自動的に認識対象語彙の切り換えが行われた場合には、何らかの方法によって話者が発声しなかった切り換え前の認識対象語彙の設定状態にもどす必要が生ずる。そして、その場合には、何らかの操作を話者に負担させるか、若しくは、自動的に切り換え前の認識対象語彙が設定されるまで話者を待たせることになるという問題がある。

【0016】そこで、この発明の目的は、自動的に認識対象語彙を切り換える場合でも高い認識精度が得られる使い易い音声認識装置および音声認識方法、並びに、音声認識処理プログラムを記録したプログラム記録媒体を提供することにある。

【0017】

【課題を解決するための手段】上記目的を達成するため、第1の発明は、入力された音声認識する認識部と、この認識部の認識結果を含む情報を出力する出力部と、上記認識時に用いられる認識対象語彙が格納された認識語彙格納部と、タイマ部と、このタイマ部からの時刻信号に基づいて上記認識対象語彙の切り換えを要求する認識対象語彙切り換え要求部を有する音声認識装置において、上記出力部は、複数の出力内容を切り換え出力するようになっており、上記認識対象語彙は、上記出力部の出力内容に対応した認識対象語の集合でなる複数の認識対象語彙セットに分類され、上記認識対象語彙の切り換えは上記認識対象語彙セットの単位で行われるようになっており、上記タイマ部からの時刻信号に基づいて、上記各認識対象語彙セット用の重みを決定する重み決定部を備えて、上記認識部は、上記全認識対象語彙セットおよび上記決定された各重みを用いて、入力音声を認識するようになっていることを特徴としている。

【0018】上記構成によれば、認識部によって、全認識対象語彙セットおよびタイマ部からの時刻信号に基づいて重み決定部によって決定された各認識対象語彙セット用の重みを用いて、入力音声認識される。その際に、上記タイマ部からの時刻信号に基づいて認識対象語彙切り換え要求部によって認識対象語彙の切り換えが要求されると、現在用いられている認識対象語彙セットが、出力部の出力内容の切り換えに応じた認識対象語彙セットに切り換えられる。したがって、切り換え前の認識対象語彙セット用の重みの値を低めるようにすれば、上記出力部の出力内容に対応している切り換え後の認識対象語彙の認識精度が高められる。

【0019】さらに、話者が、上記認識対象語彙セットが切り換えられたことを知らずに、切り換え前の認識対象語彙で発声したとしても、切り換え前の認識対象語彙セットの語をも用いて認識が行われているので、上記切り換え前の認識対象語彙セットの語に関しても高い認識精度が得られる。

【0020】また、上記第1の発明の音声認識装置は、上記重み決定部を、上記認識対象語彙切り換え要求部によって認識対象語彙の切り換えが要求されてから重み決定までの経過時間に応じて、切り換え前の認識対象語彙セット用の重みを低下させる一方、切り換え後の認識対象語彙セット用の重みを上昇させるように成すことが望ましい。

【0021】上記構成によれば、上記認識対象語彙切り換え要求部によって認識対象語彙の切り換えが要求されてからの経過時間が長くなるに連れて、切り換え前の認識対象語彙の認識精度が低くなる一方、切り換え後の認識対象語彙の認識精度が高くなる。こうして、認識に用いられる上記認識対象語彙の切り換えが徐々に行われる。

【0022】また、上記第1の発明の音声認識装置は、

上記認識部を、上記全認識対象語彙セットを構成する各語の尤度を算出し、各語の尤度の値に各語が属する認識対象語彙セット用の重みを掛け、その値が最も高い語を認識結果とするように成すことが望ましい。

【0023】上記構成によれば、認識に用いられている認識対象語彙セット用の重みと認識に用いられていない認識対象語彙セット用の重みとを最適に設定することによって、上記出力部の出力内容に対応した切り換え後の認識対象語彙の認識精度を高めることと、話者が切り換え前の認識対象語彙で発声した場合でも高い認識精度を得ることが、容易に達成される。

【0024】また、上記第1の発明の音声認識装置は、上記出力部を、上記認識対象語彙切り換え要求部からの認識対象語彙切り換え要求がなされた時点に出力している出力内容に対応する認識対象語彙セット用の重みの値と、次に出力すべき出力内容に対応する認識対象語彙セット用の重みの値との差が所定値未満になると、上記出力内容を切り換えるように成すことが望ましい。

【0025】上記構成によれば、上記認識対象語彙セットが切り換えられるのに呼応して、上記出力部の出力内容が対応する出力内容に切り換えられる。

【0026】また、第2の発明の音声認識方法は、入力された音声を認識対象語彙を用いて認識して認識結果を出力するに際して、タイマ部からの時刻信号に基づいて上記認識対象語彙の切り換えを自動的に行う音声認識方法において、複数の出力内容を出力部に切り換え出力し、上記各出力内容に対応した認識対象語の集合でなる複数の認識対象語彙セットの単位で上記認識対象語彙の切り換えを行い、上記タイマ部からの時刻信号に基づいて上記各認識対象語彙セット用の重みを決定し、上記全認識対象語彙セットおよび上記決定された各重みを用いて上記入力音声の認識を行うことを特徴としている。

【0027】上記構成によれば、全認識対象語彙セットおよびタイマ部からの時刻信号に基づいて決定された各認識対象語彙セット用の重みを用いて、入力音声認識される。その際に、上記タイマ部からの時刻信号に基づいて認識対象語彙の切り換えが要求されると、現在用いられている認識対象語彙セットが、出力部の出力内容の切り換えに応じた認識対象語彙セットに切り換えられる。したがって、切り換え前の認識対象語彙セット用の重みの値を低めるようにすれば、上記出力部の出力内容に対応している切り換え後の認識対象語彙の認識精度が高められる。

【0028】さらに、話者が、上記認識対象語彙セットが切り換えられたことを知らずに、切り換え前の認識対象語彙で発声したとしても、切り換え前の認識対象語彙セットの語をも用いて認識が行われているので、上記切り換え前の認識対象語彙セットの語に関しても高い認識精度が得られる。

【0029】また、第3の発明のプログラム記録媒体

は、コンピュータを、請求項1における認識部、出力部、タイマ部、認識対象語彙切り換え要求部および重み決定部として機能させる音声認識処理プログラムが記録されていることを特徴としている。

【0030】上記構成によれば、請求項1の場合と同様に、切り換え前の認識対象語彙セット用の重みの値を低めるようにすれば、上記出力部の出力内容に対応している切り換え後の認識対象語彙の認識精度が高められる。さらに、話者が、上記認識対象語彙セットが切り換えられたことを知らずに切り換え前の認識対象語彙で発声したとしても、高い認識精度が得られる。

【0031】

【発明の実施の形態】以下、この発明を図示の実施の形態により詳細に説明する。図1は、本実施の形態の音声認識装置におけるブロック図である。この音声認識装置21は、音声入力部22、A/D変換部23、音響分析部24、認識部25、音響モデル格納部26、第1認識対象語彙セット格納部27、第2認識対象語彙セット格納部28、重み係数決定部29、タイマ部30、認識対象語彙切り換え要求部31および出力部32で構成される。

【0032】上記音声入力部22は、マイクロホンを含む音声入力装置を備えて、入力された音声を電気信号(音声信号)に変換してA/D変換部23に出力する。A/D変換部23は、入力されたアナログ信号である音声信号をデジタル信号に変換し、デジタル化された音声信号を音響分析部24に出力する。尚、上記デジタル化された音声信号は、振幅値の時系列で表されている。

【0033】上記音響分析部24は、A/D変換部23からのデジタル音声信号からフレーム毎に特徴ベクトルを抽出して認識部25に出力する。ここで、上記特徴ベクトルは、各フレームにおける音声信号のパワー、1次～16次のLPCケプストラム係数、前フレームのパワーおよび前フレームのLPCケプストラム係数(1次～16次)の合計34の要素からなる34次元ベクトル $x_{t,i}$ を、総てのフレーム( $t=1, 2, \dots, I$ )に亘って配列したものである。

【0034】上記認識部25は、音響モデルを利用して、音響分析部24で抽出された特徴ベクトルを用いて、第1認識対象語彙セット格納部27に格納されている認識対象語彙セットAおよび第2認識対象語彙セット格納部28に格納されている認識対象語彙セットBを構成する各単語の生起確率(尤度)Pを、従来の技術で説明した手法を用いて計算する。さらに、重み係数決定部29で決定された重みwを各単語の尤度Pに掛けて、最も高い尤度 $w \cdot P$ を呈するHMMに対応する単語を出力部32に出力するのである。

【0035】上記音響モデル格納部26は、認識部25で音声認識を行う際に使用される音響モデルが格納されている。上記音響モデルは、音素を単位として、予め不特定話者の学習音声を用いてバウム・ウェルチアルゴリ

ズムと呼ばれるアルゴリズムによって学習(初期学習)されたHMMが用いられる。尚、上記HMMは、各状態における遷移確率と出力確率分布とを要素とする状態数分の配列で記憶されている。また、上記遷移確率は、各状態への遷移確率を要素として遷移数分の配列で記憶されている。また、上記出力確率は、複数の正規分布を重み付け加算した多次元の混合正規分布で表され、各正規分布における混合の重みと平均ベクトルと分散ベクトルとを要素とする次元数分の配列で記憶されている。ここで、上記平均ベクトルと分散ベクトルとは、音響分析部24で音声信号からフレーム毎に抽出される特徴ベクトルの要素数と同じ「34」の要素の配列で表される。

【0036】上記タイマ部30は、時刻を表す時刻信号を認識対象語彙切り換え要求部31、重み係数決定部29および出力部32に出力して、時刻を通知する。そうすると、認識対象語彙切り換え要求部31は、上記通知された時刻に基づいて、認識対象語彙の切り換えを要求するか否かを判断する。そして、要求する場合には、重み係数決定部29に対して認識対象語彙の切り換えを要求する。

【0037】上記重み係数決定部29は、第1認識対象語彙セット格納部27に格納されている認識対象語彙セットAおよび第2認識対象語彙セット格納部28に格納されている認識対象語彙セットBのうち、出力部32によって現在表示されている表示内容に対応する認識対象語彙セットを構成する単語に掛けられる重み $w_1$ 、および、出力部32によって表示されていない表示内容に対応する認識対象語彙セットを構成する単語に掛けられる重み $w_2$ を決定する。これらの重み $w_1, w_2$ は、記憶されている重み関数 $W_1(t), W_2(t)$ を用いて、認識対象語彙切り換え要求部31から切り換えが要求された時のタイマ部30からの時刻 $t$ を基準として所定時間 $\Delta T$ が経過する毎に決定される。そして、決定された両重み $w_1, w_2$ の値は認識部25に順次出力される。

【0038】上記第1認識対象語彙セット格納部27および第2認識対象語彙セット格納部28には、夫々の認識対象語彙セットA, Bを構成する単語が、各単語の表記と音素列との文字列を要素とする文字数分の配列で記憶されている。

【0039】上記出力部32は、ディスプレイを含む画像表示装置を備えて、認識対象語彙セットAに対応した第1表示内容と認識対象語彙セットBに対応した第2表示内容とを格納している。そして、タイマ部30から通知される時刻に基づいて、第1, 第2表示内容のうち現在表示している表示内容を変更するか否かを判断し、変更する場合は画面の表示内容を切り換える。さらに、認識部25からの認識結果を画面に表示する。

【0040】図2は、上記出力部32が現在選択している表示内容に対応する認識対象語彙セット用の重み関数 $W_1(t)$ と非選択表示内容に対応する認識対象語彙セット

用の重み関数 $W_2(t)$ との時間変化を示す。重み関数 $W_1(t)$ の値は、認識対象語彙の切り換え要求が出力された時刻 $t_0$ で1よりも小さい0近傍の所定値「a」から単調増加し始め、時刻 $t_1$ 以降は値「1」となる。一方、重み関数 $W_2(t)$ の値は、重み関数 $W_1(t)$ の値とは逆に、時刻 $t_0$ で値「1」から単調減少し始めて、時刻 $t_1$ 以降は所定値「a」となる。その場合、時刻 $t_1$ で重み $w_1$ と重み $w_2$ の差が閾値 $h$ となる。そして、出力部32は、この差の値が閾値 $h$ 未満になると、つまり認識対象語彙の切り換えが要求された時刻 $t_0$ から時間 $T(>(t_1 - t_0))$ が経過すると、画面に表示されている表示内容を切り換えるのである。

【0041】すなわち、上記出力部32がタイマ部30から通知される時刻に基づいて表示内容を変更すると判断する時点は、認識対象語彙切り換え要求部31がタイマ部30から通知された時刻に基づいて上記切り換えを要求すると判断する時点よりも上記時間 $T$ だけ遅れるように設定されているのである。

【0042】このように、本実施の形態においては、出力部32によって、自動的に画面の表示内容が切り換えられるのであるが、表示内容が切り換える前であっても切り換えた後であっても、認識部25は、認識対象語彙セットAおよび認識対象語彙セットBの両語彙セットの語彙を対象として尤度 $P$ の計算を行う。そして、現在出力部32によって選択されている表示内容に対応する認識対象語彙セットを構成する単語の尤度 $P$ には、表示内容切り換え前であれば $1 > w > (1 + a + h)/2$ であり、切り換え後であれば $1 > w > (1 + a - h)/2$ である重み $w$ を掛ける。一方、非選択側の表示内容に対応する認識対象語彙セットを構成する単語の尤度 $P$ には、表示内容切り換え前であれば $(1 + a - h)/2 > w > a$ であり、切り換え後であれば $(1 + a + h)/2 > w > a$ である重み $w$ を掛ける。こうして、最終的な尤度 $w \cdot P$ を計算して認識結果を決定するようにしている。

【0043】換言すれば、図4に示す従来の音声認識装置における認識対象語彙の切り換えは、尤度 $P$ の演算に用いる認識対象語彙そのものを切り換えることによって行うのに対して、本実施の形態においては、尤度 $P$ の演算に用いる2セットの認識対象語彙は切り換えずに尤度 $P$ に掛ける重み $w$ の値を「1」と0近傍の所定値「a」との間で徐々に変化させることによって行うのである。

【0044】したがって、本実施の形態においては、何らかの理由で話者が認識対象語彙の発声の機会を逸してしまい、且つ、自動的に認識対象語彙の切り換えが行われた後でも、切り換え前の認識対象語彙の単語に関する尤度 $w \cdot P$ の計算も行われることになり、話者が切り換え前の認識対象語彙で発声しても正しく認識することが可能になる。また、その場合、図4に示す音声認識装置のように認識対象語彙そのものを切り換えた場合と同様に、出力部32の表示内容に対応した語彙の認識精度を



高める機能は損なわれないのである。

【0045】図3は、上記重み係数決定部29によって実行される重み決定処理動作のフローチャートである。以下、図3に従って、重み決定の動作について説明する。ここで、出力部32が現在選択している表示内容に対応する認識対象語彙セット用の重み関数を $W_2(t)$ とし、非選択表示内容に対応する認識対象語彙セット用の重み関数を $W_1(t)$ とする。認識対象語彙切り換え要求部31から切り換えが要求されると重み決定処理動作がスタートする。

【0046】ステップS1で、上記タイマ部30からの時刻信号に基づいて、認識対象語彙の切り換え要求時刻 $t_0$ が取得される。ステップS2で、重み値 $w$ の算出回数 $j$ が「0」に初期化される。ステップS3で、算出回数 $j$ がインクリメントされる。ステップS4で、切り換え要求時刻 $t_0$ を取得してから又は前回重み値 $w$ を算出してから所定時間 $\Delta T$ が経過したか否かが判別される。その結果、経過していればステップS5に進む。ステップS5で、現在の時刻 $(t_0 + j \cdot \Delta T)$ が時刻 $t_1$ を越えているか否かが判別される。その結果、超えていなければステップS6に進む。

【0047】ステップS6で、上記重み関数 $W_i(t)$ の関数番号 $i$ が「1」に初期化される。ステップS7で、重み関数 $W_i(t)$ における切り換え要求時刻 $t_0$ からの経過時間 $t$ に「 $j \cdot \Delta T$ 」が代入されて、重みの値 $w_1$ が算出される。ステップS8で、関数番号 $i$ がインクリメントされる。ステップS9で、関数番号 $i$ の値が「2」よりも大きいか否かが判別される。その結果、「2」以下であればステップS7にリターンして重み値 $w_2$ の算出に移行する一方、「2」よりも大きければ、総ての認識対象語彙セットA、Bに対応する現時刻での重みが算出されたと判断されて、ステップS10に進む。ステップS10で、上記算出された現時刻での重み値 $w_1, w_2$ の配列が認識部25に出力される。そうした後、ステップS3にリターンして、次の時刻での重み値 $w_1, w_2$ の算出に移行する。

【0048】以後、上記ステップS3～ステップS10を繰り返し、ステップS5において現在の時刻が時刻 $t_1$ を越えていると判別されると、重み決定処理動作を終了する。その後は、表示内容に対応する認識対象語彙セット用の重み値 $w_2$ として「1」が所定時間 $\Delta T$ 毎に出力され、非選択表示内容に対応する認識対象語彙セット用の重み値 $w_1$ として所定値「a」が所定時間 $\Delta T$ 毎に出力される。そして、次に認識対象語彙切り換え要求部31から切り換え要求が出力されると、上記重み決定処理動作がスタートするのである。

【0049】上述のように、本実施の形態における認識部25は、音響モデル格納部26に格納された音響モデルを用いて、第1認識対象語彙セット格納部27に格納された認識対象語彙セットAと第2認識対象語彙セット

格納部28に格納された認識対象語彙セットBとを構成する単語の尤度 $P$ を算出する。その際における出力部32の表示内容の切り換えに伴う認識対象語彙セットの切り換えは、認識対象語彙セットそのものを切り換えるのではなく、選択、非選択認識対象語彙セットを構成する単語の尤度 $P$ に掛ける重み $w_2, w_1$ の値を「1」と0近傍の所定値「a」とに切り換えることによって行う。そして、その場合に、重み $w_2, w_1$ の値を段階的に切り換えるのではなく、認識対象語彙切り換え要求部31から切り換え要求がなされた時刻 $t_0$ からの経過時間「 $j \cdot \Delta T$ 」に比例して徐々に値「1」から値「a」へ又は値「a」から値「1」へ切り換えるようにしている。

【0050】したがって、本実施の形態によれば、何らかの理由で話者が認識対象語彙の発声の機会を逸してしまい、且つ、自動的に認識対象語彙が切り換えられてしまっても、切り換え前の認識対象語彙セットの単語に関する尤度 $w \cdot P$ の計算をも行うので、話者が切り換え前の認識対象語彙で発声しても正しく認識することができる。また、その場合に、図4に示す音声認識装置のごとく認識対象語彙そのものを切り換える場合と同様に、出力部32の表示内容に対応した認識対象語彙の認識精度を高める機能は損なわれることはない。

【0051】尚、上記実施の形態においては、選択認識対象語彙セット用の重み関数 $W_2(t)$ および非選択表示内容に対応する認識対象語彙セット用の重み関数 $W_1(t)$ を、認識対象語彙切り換え要求部31による切り換え要求時刻 $t_0$ からの経過時間「 $j \cdot \Delta T$ 」に比例して、値「1」、「a」から値「a」、「1」へ直線的に切り換えるようにしている。しかしながら、この発明においては、関数 $W_1(t), W_2(t)$ の形状は直線に限定されるものではない。曲線にして、表示内容の切り換え時刻 $t_1$ までの関数 $W_2(t)$ の値を高める一方関数 $W_1(t)$ の値を低め、表示内容の切り換え時刻 $t_1$ 以降の関数 $W_2(t)$ の値を低める一方関数 $W_1(t)$ の値を高めてもよい。

【0052】また、上記実施の形態においては、上記重み係数決定部29を、認識対象語彙切り換え要求部31からの切り換え要求時刻 $t_0$ を基準として所定時間 $\Delta T$ が経過する毎に重み値 $w_1, w_2$ を決定して認識部25に出力するように構成し、認識部25は、入力される重み値 $w_1, w_2$ を必要に応じて用いて認識処理を行うように構成している。しかしながら、この発明はこれに限定されるものではなく、認識部25を、認識を行う際に重み係数決定部29に対して重み決定要求を出すように構成し、重み係数決定部29は、重み決定要求を受けると、認識対象語彙切り換え要求部31による切り換え要求時刻 $t_0$ からの経過時間を重み関数 $W_i(t)$ に代入して算出するように構成しても差し支えない。

【0053】ところで、上記各実施の形態における上記認識部、出力部、タイマ部、認識対象語彙切り換え要求部および重み決定部としての機能は、プログラム記録媒体

10

20

30

40

50

に記録された音声認識処理プログラムによって実現される。上記実施の形態における上記プログラム記録媒体は、ROM(リード・オンリ・メモリ)でなるプログラムメディアである。あるいは、外部補助記憶装置に装着されて読み出されるプログラムメディアであってもよい。尚、何れの場合においても、上記プログラムメディアから音声認識処理プログラムを読み出すプログラム読み出し手段は、上記プログラムメディアに直接アクセスして読み出す構成を有していてもよいし、RAM(ランダム・アクセス・メモリ)に設けられたプログラム記憶エリア(図示せず)にダウンロードし、上記プログラム記憶エリアにアクセスして読み出す構成を有していてもよい。尚、上記プログラムメディアからRAMの上記プログラム記憶エリアにダウンロードするためのダウンロードプログラムは、予め本体装置に格納されているものとする。

【0054】ここで、上記プログラムメディアとは、本体側と分離可能に構成され、磁気テープやカセットテープ等のテープ系、フロッピー(登録商標)ディスク、ハードディスク等の磁気ディスクやCD(コンパクトディスク)・ROM,MO(光磁気)ディスク,MD(ミニディスク),DVD(デジタルビデオディスク)等の光ディスクのディスク系、IC(集積回路)カードや光カード等のカード系、マスクROM,EPROM(紫外線消去型ROM),EEPROM(電氣的消去型ROM),フラッシュROM等の半導体メモリ系を含めた、固定的にプログラムを担持する媒体である。

【0055】また、上記各実施の形態における音声認識装置は、モデムを備えてインターネットを含む通信ネットワークと接続可能な構成を有していれば、上記プログラムメディアは、通信ネットワークからのダウンロード等によって流動的にプログラムを担持する媒体であっても差し支えない。尚、その場合における上記通信ネットワークからダウンロードするためのダウンロードプログラムは、予め本体装置に格納されているものとする。あるいは、別の記録媒体からインストールされるものとする。

【0056】尚、上記記録媒体に記録されるものはプログラムのみに限定されるものではなく、データも記録することが可能である。

【0057】

【発明の効果】以上より明らかなように、第1の発明の音声認識装置は、出力部の出力内容に対応した複数の認識対象語彙セットを認識語彙格納部に格納し、重み決定部によって、タイマ部からの時刻信号に基づいて上記各認識対象語彙セット用の重みを決定し、認識部によって、上記全認識対象語彙セットおよび上記決定された各重みを用いて入力音声を認識するので、認識対象語彙切り換え要求部による認識対象語彙の切り換え要求に基づいて、上記出力部の出力内容の切り換えに応じた認識対

象語彙セットに切り換えられる際に、切り換え前の認識対象語彙セット用の重みの値を低めるようにすれば、切り換え後の認識対象語彙の認識精度を高めることができる。

【0058】さらに、話者が、上記認識対象語彙セットが切り換えられたことを知らずに、切り換え前の認識対象語彙で発声しても、切り換え前の認識対象語彙セットの語をも用いて認識を行うので、上記切り換え前の認識対象語彙セットの語に関しても高い認識精度を得ることができる。

【0059】すなわち、この発明によれば、自動的に認識対象語彙を切り換える場合でも高い認識精度を得ることができる。さらに、その際に話者に何らかの操作や待ち時間を負担させることがなく、使い易い音声認識装置を実現できる。

【0060】また、上記第1の発明の音声認識装置は、上記重み決定部を、上記認識対象語彙切り換え要求部によって認識対象語彙の切り換えが要求されてから重み決定までの経過時間に応じて、切り換え前の認識対象語彙セット用の重みを低下させる一方、切り換え後の認識対象語彙セット用の重みを上昇させるように成せば、認識に用いられる上記認識対象語彙の切り換えを徐々に行うことができる。したがって、上記切り換え前の認識対象語彙セットの語に関しても高い認識精度を得ることができる。

【0061】また、上記第1の発明の音声認識装置は、上記認識部を、全認識対象語彙セットを構成する各語の尤度を算出し、各語の尤度の値に各語が属する認識対象語彙セット用の重みを掛け、その値が最も高い語を認識結果とするように成せば、認識に用いられている認識対象語彙セット用の重みと認識に用いられていない認識対象語彙セット用の重みとを最適に設定すれば、上記出力部の出力内容に対応している切り換え後の認識対象語彙の認識精度を高めることと、話者が切り換え前の認識対象語彙で発声した場合でも高い認識精度を得ることとを、容易に達成することができる。

【0062】また、上記第1の発明の音声認識装置は、上記出力部を、上記認識対象語彙切り換え要求部からの認識対象語彙切り換え要求がなされた時点に出力している出力内容に対応する認識対象語彙セット用の重みの値と、次に出力すべき出力内容に対応する認識対象語彙セット用の重みの値との差が所定値未満になると、上記出力内容を切り換えるように成せば、上記認識対象語彙セットが切り換えられるのに呼応して、上記出力部の出力内容に対応する出力内容に切り換えることができる。

【0063】また、第2の発明の音声認識方法は、タイマ部からの時刻信号に基づいて出力部の出力内容に対応した複数の認識対象語彙セット用の重みを決定し、全認識対象語彙セットおよび上記決定された各重みを用いて入力音声を認識するので、認識対象語彙セットが切り換

えられる際に、切り換え前の認識対象語彙セット用の重みの値を低めるようにすれば、上記出力部の出力内容に応じた切り換え後の認識対象語彙の認識精度を高めることができる。

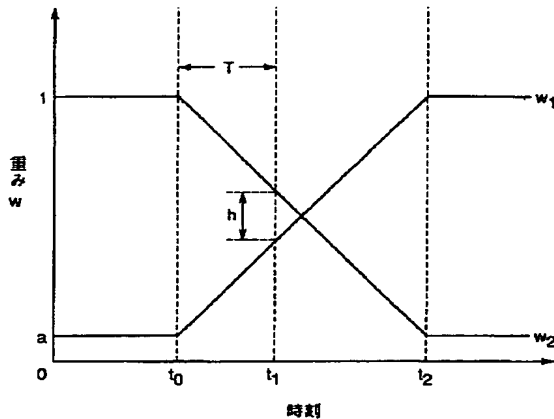
【0064】さらに、話者が、上記認識対象語彙セットが切り換えられたことを知らずに、切り換え前の認識対象語彙で発声しても、切り換え前の認識対象語彙の語をも用いて認識を行うので、上記切り換え前の認識対象語彙セットの語に関しても高い認識精度を得ることができる。

【0065】また、第3の発明のプログラム記録媒体は、コンピュータを、請求項1における認識部、出力部、タイマ部、認識対象語彙切り換え要求部および重み決定部として機能させる音声認識処理プログラムが記録されているので、請求項1の場合と同様に、切り換え前の認識対象語彙セット用の重みの値を低めるようにすれば、上記出力部の出力内容に対応している切り換え後の認識対象語彙の認識精度を高めることができる。さらに、話者が、上記認識対象語彙セットが切り換えられたことを知らずに切り換え前の認識対象語彙で発声したとしても、高い認識精度を得ることができる。

【図面の簡単な説明】

\*

【図2】



\*【図1】 この発明の音声認識装置におけるブロック図である。

【図2】 選択、非選択認識対象語彙セット用の重み関数の時間変化を示す図である。

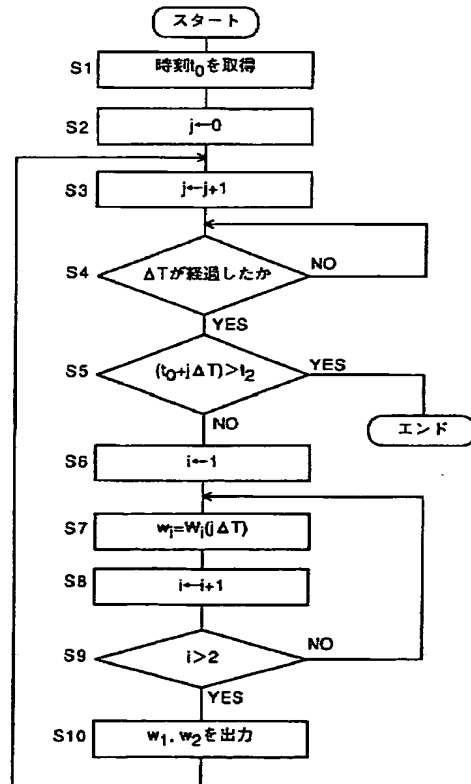
【図3】 図1における重み係数決定部によって実行される重み決定処理動作のフローチャートである。

【図4】 認識対象語彙の切り換えが可能な従来の音声認識装置のブロック図である。

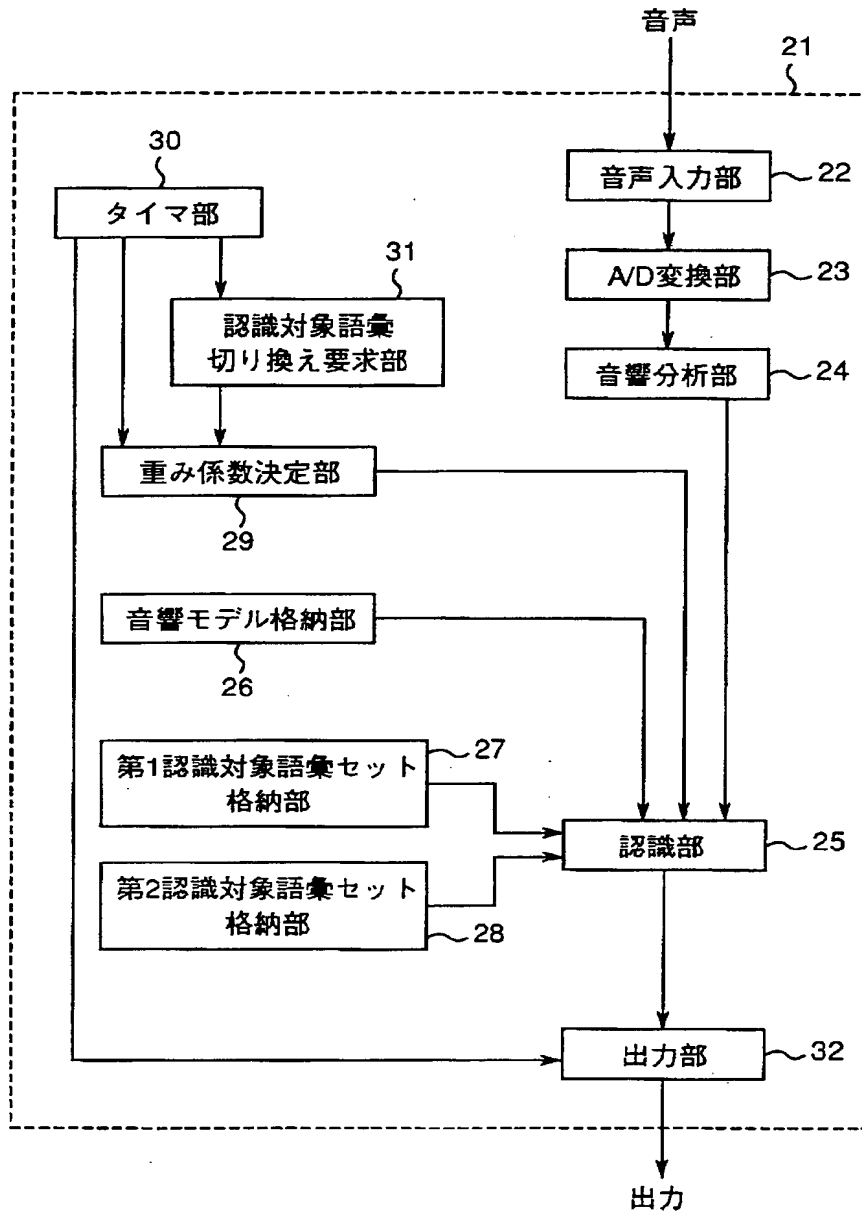
【符号の説明】

- 10 21…音声認識装置、
- 22…音声入力部、
- 23…A/D変換部、
- 24…音響分析部、
- 25…認識部、
- 26…音響モデル格納部、
- 27…第1認識対象語彙セット格納部、
- 28…第2認識対象語彙セット格納部、
- 29…重み係数決定部、
- 30…タイマ部、
- 31…認識対象語彙切り換え要求部、
- 32…出力部。

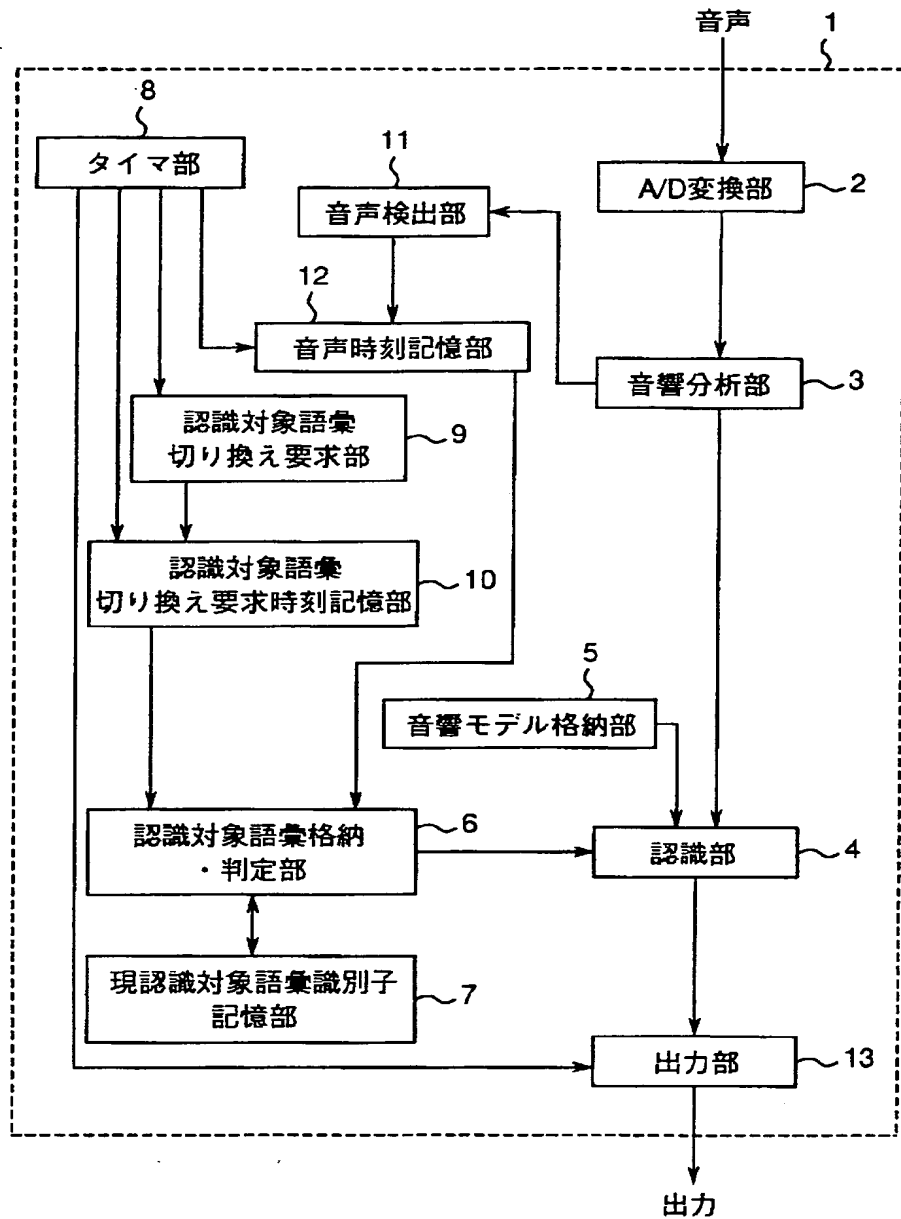
【図3】



【図1】



【図4】



フロントページの続き

(72)発明者 勘座 浩幸  
 大阪府大阪市阿倍野区長池町22番22号 シ  
 ャープ株式会社内

Fターム(参考) 5D015 AA04 BB02 HH05 HH11 HH21  
 KK02 LL10

This Page Blank (uspto)